

CLAIMS

What is claimed is:

1. A method for identifying a repeat sequence, the method comprising the steps of:
 - selecting a query sequence;
 - testing said query sequence with a redundant file;
 - identifying sequences in the redundant file that contain a similar sequence to a portion of the query sequence, wherein said identified sequences and said similar portion of the query sequence make up a pairwise sequence alignment;
 - aligning all the identified pairwise sequence alignments;
 - designating the right and left endpoints of each identified sequence and any intervening sequences;
 - identifying a position within the query sequence corresponding to each endpoint;
 - defining regions within the query sequence, wherein a region is a sequence between two consecutive positions matching two endpoints; and
 - identifying each regions having at least five sequence matches in the identified pairwise alignments as a repeat sequence.
2. A method for constructing a repeat database comprising:
 - selecting a query sequence;
 - selecting known repeat sequences;
 - adding known repeat sequences into a repeat sequence database;
 - masking said query sequence with repeat sequences in the repeat sequence database;
 - testing said masked query sequence with a redundant file;

1 identifying sequences in the redundant file that contain a similar sequence to a portion of
2 the query sequence, wherein said identified sequences and said similar portion of the query
3 sequence make up a pairwise sequence alignment;

4 aligning all the identified pairwise sequence alignments;

5 designating the right and left endpoints of each identified sequence and any intervening
6 sequences;

7 identifying a position within the query sequence corresponding to each endpoint;

8 defining regions within the query sequence, wherein a region is a sequence between two
9 consecutive positions matching two endpoints;

10 identifying any two successive regions having a large variance in the number of sequence
11 matches; and

12 adding the sequence within the region of the two successive regions having the highest
13 number of sequence matches into the repeat sequence database.

14
15 3. The method of claim 2, wherein the large variance in the number of sequence matches is
16 equal to 5 or more.

17
18 4. A database product of the process of claim 2.

19
20 5. The method of claim 1 or 2, wherein said sequence is a deoxyribonucleotide sequence.

21
22 6. The method of claim 1 or 2, wherein said sequence is a ribonucleotide sequence.

1 7. The method of claim 1 or 2, wherein said sequences are derived from animal DNA or
2 RNA.

4 8. The method of claim 7, wherein said animal is a human.

6 9. The method of claim 8, wherein said animal is a mouse.

8 10. The method of claim 1 or 2, wherein said sequences are derived from plant DNA or
9 RNA.

11 11. The method of claim 10, wherein said plant is a single-cell plant.

13 12. The method of claim 1 or 2, wherein said sequences are derived from fungal DNA or
14 RNA.

16 13. The method of claim 1 or 2, wherein said sequences are derived from DNA or RNA of a
17 microorganism or virus.

19 14. The method of claim 1 or 2, wherein said sequences are derived from DNA or RNA of a
20 single-cell eukaryote.

22 15. The method of claim 1 or 2, wherein said sequences are derived from synthetic man-
23 made DNA or RNA.

1 16. The method of claim 1 or 2, wherein said sequences are postulated based upon amino
2 acid sequences.

3
4 17. The method of claim 2, wherein said database is encoded in a biological medium.

5
6 18. The method of claim 2, wherein said database is encoded in a written medium.

7
8 19. The method of claim 2, wherein said database is encoded in an electronic medium.

9
10 20. The method of claim 19, wherein said electronic medium is a computer-readable
11 medium.

12
13 21. The method of claim 20, wherein said computer-readable medium is addressable through
14 an internet connection.

15
16 22. The method of claim 1 or 2, wherein said redundant file is a Public Domain Database.

17
18 23. The method of claim 22, wherein said Public Domain Database is GenBank.

19
20 24. The method of claim 22, wherein said Public Domain Database is dbEST.

21
22 25. The method of claim 22, wherein said Public Domain Database is TIGR.

23
24 26. The method of claim 22, wherein said Public Domain Database is SwissProt.

1 27. The method of claim 1 or 2, wherein sequence comparisons are carried out using a
2 Database Search Algorithm.

4 28. The method of claim 27, wherein said Database Search Algorithm is BLAST.

6 29. The method of claim 27, wherein said Database Search Algorithm is FASTA.

8 30. The method of claim 27, wherein said Database Search Algorithm is Smith-Waterman.

10 31. The method of claim 1 or 2, wherein said sequence comparisons are carried out utilizing
a Scoring Matrix Program.

13 32. The method of claim 31, wherein said Scoring Matrix Program is PAM.

15 33. The method of claim 31, wherein said Scoring Matrix Program is BLOSUM.

17 34. The process of Figure 2.

19 35. A repeat sequence product of the process of claim 1.

21 36. A kit for analyzing nucleotide sequences comprising:

22 an electronic medium readable by a computer, said medium encoding a database
23 produced by the method of claim 2.

25 37. A kit for analyzing nucleotide sequences comprising:

1 an electronic medium readable by a computer, said medium encoding a database
2 produced by the method of claim 2; and,
3 instructions for the use of said database.
4

5 38. A kit for analyzing nucleotide sequences comprising:

6 an electronic medium readable by a computer, said medium encoding a database
7 produced by the method of claim 2;
8 instructions for the use of said database; and,
9 a computer.
10

11 39. An improved database of nucleotide sequences, the improvement consisting of repeat
12 sequences containing a similar sequence to a portion of a query sequence, wherein said identified
13 sequences and said similar portion of the query sequence make up a pairwise sequence
14 alignment, and wherein all identified pairwise sequence alignments have right and left endpoints
15 of each identified sequence and any intervening sequences.